

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2001-051890

(43)Date of publication of application : 23.02.2001

(51)Int.Cl.

G06F 12/00

G06F 13/00

(21)Application number : 11-226494

(71)Applicant : TOSHIBA CORP

(22)Date of filing : 10.08.1999

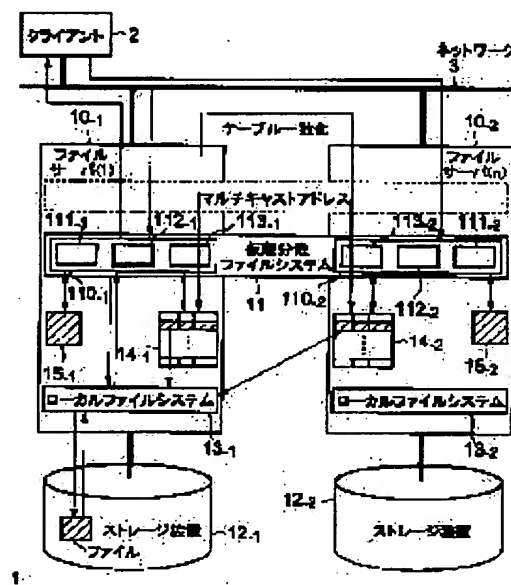
(72)Inventor : UCHIBORI IKUO
TAKAKUWA MASAYUKI

(54) VIRTUAL DECENTRALIZED FILE SERVER SYSTEM

(57)Abstract:

PROBLEM TO BE SOLVED: To make a client not pay attention to the number of file servers decentralized in a network and the connection states of storage devices.

SOLUTION: This virtual decentralized file server system 1 is equipped with servers 10-1 and 10-2 decentralized in the network 3 and a virtual decentralized file system 11 is decentralized and mounted on each of the servers. Modules 110-1 and 110-2 on the servers 10-1 and 10-2 which constitute this system 11 when receiving a file operation request multicast from a client 2 judge whether or not their servers are optimum servers capable of handling the request according to server information holding parts 15-1 and 15-2 holding mapping tables 14-1 and 14-2 between the virtual decentralized file system 11 and all local file systems 13-1 and 13-2 or server information on all the servers, and makes a local file system of a corresponding server perform requested file operation according to the judgement result.



LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of rejection]

[Date of requesting appeal against examiner's decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

Japanese Patent Laid-open Publication No. 2001-51890 A

Publication date: February 23, 2001

Applicant: TOSHIBA CORP

Title: Virtual distributed file server system

5

(57) [Abstract]

[Object] To make it unnecessary that clients consider the number of file servers distributed on a network and the connection state of storage devices.

[Means] Distributed servers 10-1 and 10-2 are provided on a network 3, and a virtual distributed file system 11 is distributed and mounted on each server. When modules 110-1 and 110-2 on the servers 10-1 and 10-2 constituting this system 11 receive a multi-cast file operation request from a client 2, determine whether the own server is the optimum server capable of processing the request, based on mapping tables 14-1, 14-2 between the virtual distributed file system 11 and all local file systems 13-1 and 13-2, or server information holding sections 15-1, 15-2 that hold server information of all servers, and allow a local file system of the corresponding server to perform the requested file operation based on the determination result.

[What is claimed is]

20 [Claim 1] A virtual distributed file server system comprising a plurality of file servers capable of multi-casting and distributed on a network, comprising:

a virtual distributed file system distributed and mounted on the respective file servers and independent of the actual storage configuration, which manages files in all file servers in an integrated manner;

25 a local file system mounted on each file server respectively independently, for

managing the storage configuration peculiar to the respective servers; and

a mapping table provided in the respective file servers for holding mapping information between the virtual distributed file server system and the local file system that actually manages files, with respect to respective files managed by the virtual

5 distributed file system in an integrated manner;

wherein the virtual distributed file system comprises management modules respectively provided in the respective file servers, and each management module

commonly receives a multi-cast file operation request from a client, and refers to the mapping table in the own server in response to the request, to determine whether the

10 own server is the optimum server capable of processing the request, and only when it determines that the own server is the optimum server, allows the local file system in the corresponding server to perform the requested file operation.

[Claim 2] A virtual distributed file server system comprising a plurality of file servers capable of multi-casting and distributed on a network, comprising:

15 a virtual distributed file system distributed and mounted on the respective file servers and independent of the actual storage configuration, which manages files in all file servers in an integrated manner;

a local file system mounted on each file server respectively independently, for managing the storage configuration peculiar to the respective servers;

20 a mapping table provided in the respective file servers for holding mapping information between the virtual distributed file server system and the local file system that actually manages files, with respect to respective files managed by the virtual distributed file system in an integrated manner; and

a server information holding unit provided in the respective file servers, which
25 holds, with respect to all file servers, server information including at least one of

information indicating unused capacity of a storage device in the server, and
information indicating the load status of the server,

wherein the virtual distributed file system comprises management modules
respectively provided in the respective file servers, and each management module
5 commonly receives a multi-cast file operation request from a client, and refers to the
mapping table in the own server or the server information holding unit in response to
the request, to determine whether the own server is the optimum server capable of
processing the request, and only when it determines that the own server is the
optimum server, allows the local file system in the corresponding server to perform the
10 requested file operation.

[Claim 3] The virtual distributed file server system according to claim 1 or 2,
wherein when the file operation request is a file read request or a file write request, the
management module refers to the mapping table in the own server, to determine
whether the own server is the optimum server capable of processing the request,
15 according to whether the corresponding file is under management of the local file
system in the own server.

[Claim 4] The virtual distributed file server system according to claim 2, wherein
when the file operation request is a new file creation request, the management module
refers to the server information holding unit in the own server, to compare the unused
20 capacity of the storage device of the server or the load status of the server, for each of
all servers, thereby determining whether the own server is the optimum server capable
of processing the request.

[Claim 5] The virtual distributed file server system according to claim 1, wherein
the management modules exchange the information of the mapping table in the own
25 server and the information of the mapping table in the other servers by inter-server

communication, in order to make the contents of the mapping tables in all file servers agree with each other.

[Claim 6] The virtual distributed file server system according to claim 2, wherein the management modules exchange the information of the mapping table in the own
5 server and the information of the mapping table in the other servers by inter-server communication, in order to make the contents of the mapping tables in all file servers agree with each other, and exchange the information of the server information holding unit in the own server and the information of the server information holding unit in the other servers by inter-server communication, in order to make the contents of the
10 server information holding unit in all file servers agree with each other.

[Claim 7] The virtual distributed file server system according to claim 1 or 2, further comprising an each file load status information holding unit respectively provided in the respective file servers for holding the information indicating the load status for each file under management of the server, wherein

15 the management module detects a file having a load exceeding a first threshold from information held by the each file load status information holding unit in the own server, to replicate the file with respect to another optional file server by inter-server communication, and when there is a multi-cast read request of the file from a client, entrusts the replicated side with the processing with respect to the request.

20 [0004]

[Problems to be solved by the Invention]

In the computer network system, when files are shared by using a file server, it is general that the client side can see the volume configuration on the file server side.

25 For example, when the server side increases a disk drive, the client side must

recognize the new volume and mount the volume. When the server itself is increased, complicated operation such as determination of the operation policy of the increased server or system setup and management becomes necessary, and the client side must recognize the new server and mount the new volume.

5 [0005] In the file sharing system using the conventional file server (file server system), when an increase of a disk drive (storage device) or an increase of a server is required, there is a problem in that great cost is required for new setup and management both on the server and client sides. Further, according to the operation mode of the storage, it may be desired to expand only the capacity of a specific file system, or only
10 an increase of the storage device and the server may not solve the problem.

[0006] The present invention has been achieved in order to solve the above problems. It is an object of the present invention to provide a virtual distributed file server system in which clients can use a plurality of file servers distributed on a network as a single server, and clients are unnecessary to consider the number of servers and
15 the connection state of the storage devices.

[0028]

[First Embodiment]

Fig. 1 is a block diagram illustrating the configuration of a computer network
20 system to which the virtual distributed file server system according to the first embodiment of the present invention is applied.

[0029] In this figure, reference sign 1 denotes a virtual distributed file server system, and 2 denotes a client (a client computer) that requests file service to the file server system 1. The virtual distributed file server system 1 is realized by using a plurality of,
25 for example, two file servers (server computers) 10-1, 10-2 distributed on a network 3.

In this figure, for the convenience sake, only one client 2 is shown, but it is general that a plurality of clients exists.

[0030] Reference sign 11 denotes a virtual distributed file system, being the center of the virtual distributed file server system 1, which is distributed and mounted on the
5 respective file servers 10-1 and 10-2. The virtual distributed file system 11 manages files in all file servers 10-1 and 10-2 in an integrated manner, and provides the client 2 with a virtual file system that does not depend on the actual volume configuration (storage configuration) of the respective file servers 10-1 and 10-2.

[0031] The virtual distributed file system 11 has virtual distributed file modules 110-1
10 and 110-2 distributed and mounted on the respective file servers 10-1 and 10-2. The virtual distributed file modules 110-1 and 110-2 are management modules for virtually making the file servers look like one file system with respect to the client 2, while distributing and processing the request from the client 2 on the file servers 10-1 and 10-2. The virtual distributed file modules 110-1 and 110-2 have virtual distributed file
15 interfaces 111-1 and 111-2, which is the center of the modules 110-1 and 110-2, and processes the request from the client 2, interfaces (local file interfaces) 112-1 and 112-2 between local file systems 13-1 and 13-2 and the modules, and communication modules 113-1 and 113-2 that performs communication between the modules 110-1 and 110-2 (communication between servers represented by communication for
20 making the information of the mapping tables 14-1, 14-2 and the server information agree with each other).

[0032] The file servers 10-1 and 10-2 are connected to the client 2 via the network 3. The file servers 10-1 and 10-2 are mounted with local file systems 13-1 and 13-2 that manage the storage devices 12-1 and 12-2 (actual storage configuration), such as disk
25 drives respectively connected to the servers 10-1 to 10-n, in addition to the virtual

distributed file system 11.

[0033] The file servers 10-1 and 10-2 are provided with the mapping tables 14-1 and 14-2 having the same content for associating the virtual distributed file server system 111 with the local file systems 13-1 and 13-2. The data structure of the table 14-i ($i = 1, 2$) is shown in Fig. 2.

[0034] Each entry to the table 14-i includes a file name registration field (file name field) 141 of a file managed by the virtual distributed file system 1, a registration field (a virtual path field) 142 of a path (a virtual path) expressing a logical position (which can be seen from the client 2) of the file on the virtual distributed file system 1, a registration field (position information field) 143 of position information expressing a physical position (which cannot be seen from the client 2) of the file on the storage, a registration field (permission information field) 144 of permission information for managing the access right (permission/prohibition) to the file, and a registration field 145 of other various attributes.

[0035] The (virtual distributed file module 110-i on the) virtual distributed file system 11 refers to the mapping table 14-i having such a data structure, thereby obtaining location information, for example, in which file server 10-1 or 10-2 a certain file exists, and also obtain the attribute of the file, such as permission, according to need.

[0036] The file servers 10-1 and 10-2 are provided with server information holding sections 15-1 and 15-2 having the same content. The server information holding section 15-i ($i = 1, 2$) is used, as shown in Fig. 3, for holding server information including information indicating the unused storage capacity (resource information) of (the storage devices 12-1, 12-2 of) all file servers 10-1 and 10-2 constituting the virtual distributed file server system 1, and information indicating the load status.

[0037] The operation of the configuration shown in Fig. 1 will be explained next. In

the embodiment, it seems to the client 2 that not the local file systems 13-1 and 13-2 of the respective file servers 10-1 and 10-2, but the virtual distributed file system 11 is mounted. Therefore, when any file operation request occurs, the client 2 issues the same request to all file servers 10-1 and 10-2 on which the virtual distributed file
5 system 11 is mounted. In this case, for example, according to a method such as using Internet protocol (IP) multi-cast, the client 2 side can issue a request without being aware of the number of the file servers.

[0038] Upon reception of the request from the client 2, the file servers 10-1 and 10-2 hand over the request to virtual distributed file modules 110-1 and 110-2 corresponding
10 to the own server in the virtual distributed file system 11. (The virtual distributed file interfaces 111-1 and 111-2 in) the modules 110-1 and 110-2 determine the request type, whether the request is a file read request or a file write (update) request, or a new file creation request or a directory creation request.

[0039] It is assumed herein that the file operation request from the client 2 is a file
15 read request or a file write request. The request is added with a file name of the requested file, and a path (virtual path) of the file on the virtual distributed file system 11.

[0040] When the file operation request from the client 2 is a file read request or a file write request, (the virtual distributed file interfaces 111-1 and 111-2 in) the virtual distributed file modules 110-1 and 110-2 refer to the mapping tables 14-1, 14-2 in the
20 own server, by the file name and the virtual path of the requested file, to check if the operation-requested file is held in the own server (in the storage devices 12-1, 12-2 connected to the own server), based on the file name and the registration information of the position information field 143 in the entry of the tables 14-1 and 14-2.

[0041] If the requested file is held in the own server, (the virtual distributed file
25 interfaces 111-1, 111-2 in) the virtual distributed file modules 110-1, 110-2 access the

actual file via the local file systems 13-1, 13-2 in the own server by the local file interfaces 111-1, 111-2, and give a response to the client 2. On the other hand, if the operation-requested file is not held in the own server, the virtual distributed file modules 110-1, 110-2 regard that another server is to respond, and do not give a response.

- 5 [0042] On the other hand, if the request from the client 2 is a new file creation or a directory creation, the virtual distributed file modules 110-1, 110-2 refer to the server information holding sections 15-1, 15-2 (not to the mapping tables 14-1, 14-2).

According to a predetermined algorithm, based on the server information of all servers held in the server information holding sections 15-1, 15-2, only the virtual distributed file
10 modules 110-i on either one server 10-i (i is 1 or 2) accepts the request from the client 2. Specifically, the virtual distributed file modules 110-i on the server 10-i compares the unused storage capacities indicated by the server information of all servers, and when it can be determined that the unused storage capacity of (the storage device 12-i of) the own server 10-i is the largest, accepts the request from the client 2. In this case, it
15 is not always necessary that the information of the unused storage capacity of the corresponding server is included in the server information.

[0043] The virtual distributed file module 110-i may compare the loads indicated by the server information of all servers, and when the load of the own server is the lowest, may accept the request from the client 2. In this case, it is not always necessary that
20 the information of the unused storage capacity of the corresponding server is included in the server information.

[0044] Other than this configuration, the virtual distributed file module 110-i may determine continuous areas that can be secured on the respective storage devices 12-1, 12-2, from the mapping information for each file in the mapping table 14-i (that is,
25 from the status of use of the area in the respective storage devices 12-1, 12-2), and

when continuous areas not smaller than the necessary size can be secured, and the storage device having the largest size is the storage device 12-i in the own server, may accept the request from the client 2. In this case, the server information holding sections 15-1, 15-2 may not be always necessary.

- 5 [0045] Further, an evaluation value may be obtained under the condition of at least two of the unused storage capacity, the load status, and the size of the securable continuous areas (complex condition), to determine whether the own server is the optimum server for accepting the request.

[0046] In the virtual distributed file module 110-i on the file server 10-i, upon reception
10 of the request from the client 2 by the virtual distributed file interface 111-i, the local file system 13-i creates a requested new file or a directory via the local file interface 112-i, thereby registering the corresponding entry information in the mapping tables 14-1 and 14-2.

[0047] After the new file or the directory has been created, the virtual distributed file
15 module 110-i on the file server 10-i transmits the new entry information registered in the mapping table 14-i on the own server to the virtual distributed file module 110-j (j is 1 or 2, provided that $j \neq i$) on all other servers 10-j, by the communication module 113-i via the network 3. (The virtual distributed file interface 111-j in) the virtual distributed file module 110-j receives the entry information for the mapping table 14-i transmitted from
20 the virtual distributed file module 110-i, and registers the entry information for the mapping table 14-i of the received other server 10-i, in the mapping table 14-j in the own server. Thus, since the virtual distributed file modules 110-1 and 110-2 on the file servers 10-1 and 110-2 exchange the newly registered (and updated) entry information for the mapping tables 14-1 and 14-2 with each other, the contents in the mapping
25 tables 14-1 and 14-2 can be made to agree with each other.

[0048] The virtual distributed file modules 110-1 and 110-2 on the file servers 10-1 and 10-2 regularly update the server information (unused storage capacity and load status) of the own server, of the server information of the own server held by the server information holding sections 15-1 and 15-2 on the own server, and regularly transmit
5 the updated server information to (the virtual distributed file modules 110-1 and 110-2 on) all other servers, by the communication modules 113-1, 113-2) via the network 3, thereby making the contents of the server information holding sections 15-1 and 15-2 in the respective file servers 10-1 and 10-2 agree with each other. In other words, the virtual distributed file modules 110-1 and 110-2 make the contents agree with each
10 other by exchanging the server information regularly.

[0049] By the above operation, the file servers 10-1 and 10-2 can be autonomously distributed and made to cooperate with each other. As a result, the client 2 can be provided with a virtual file server, without noticing that two (a plurality of) file servers actually exist.

15 [0050] In the system example in Fig. 1, a system having two servers has been explained, but even when the system includes three or more servers, the virtual file server can be also provided by the similar mechanism.

[0051]

[Second Embodiment]

20 Fig. 4 is a block diagram illustrating the configuration of a computer network system to which the virtual distributed file server system according to the second embodiment of the present invention is applied, wherein like reference signs refer to like parts in Fig. 1.

[0052] In Fig. 4, a virtual distributed file system 41, being the center of a virtual
25 distributed file server system 4, is distributed and mounted on n file servers 10-1 to

10-n. As the virtual distributed file system 11 in Fig. 1, the virtual distributed file system 41 manages files in all file servers 10-1 to 10-n in an integrated manner, and provides the client 2 with a virtual file system that does not depend on the actual volume configuration of the respective file servers 10-1 to 10-n. The virtual distributed file system 41 has virtual distributed file modules 410-1 and 410-n that process the request from the client 2 on the respective file servers 10-1 to 10-n. The modules 410-1 and 410-n have the same configuration as the modules 110-1 and 110-2 in Fig. 1, and includes virtual distributed file interfaces 411-1 to 411-n, local file interfaces 412-1 to 412-n, and communication modules 413-1 to 413-n. However, the communication modules 413-1 to 413-n in this embodiment are different from the communication modules 113-1 and 113-2 in Fig. 1, and constructed so as to perform communication via a private communication path 5 described later.

[0053] The file servers 10-1 to 10-n are connected with the client 2 via the network 3. The file servers 10-1 to 10-n are mounted with local file systems (local file systems) 13-1 to 13-2 that manage the storage devices 12-1 to 12-2 respectively connected to the corresponding servers 10-1 to 10-n. The file servers 10-1 to 10-n are provided with the mapping tables 14-1 to 14-n and the server information holding sections 15-1 to 15-2.

[0054] The characteristic points of the virtual distributed file server system 4 having the configuration shown in Fig. 4 are that the number of file servers constituting the system is n, and that the n file servers 10-1 to 10-n are connected with each other by the private communication path 5 separately from the network 3, different from the virtual distributed file server system 1 shown in Fig. 1. The private communication path 5 is, for example, Ethernet or Fibre Channel, but the physical layer is not particularly specified. As for the topology, the bus type topology is assumed in the

example in Fig. 4, but it may be a loop or a switch.

[0055] In the configuration in Fig. 4, in order that the file servers 10-1 to 10-n perform distribution and cooperation operation (by the virtual distributed file modules 410-1 to 410-n in the virtual distributed file system 41), it is necessary to make the contents of the mapping tables 14-1 to 14-n and the server information holding sections 15-1 to 15-n agree with each other between respective servers 10-1 to 10-n at all times, as can be figured out from the explanation in the first embodiment. However, if the agreement of information between the servers 10-1 to 10-n is realized via the network 3 as in the first embodiment, when the number of file servers constituting the virtual distributed file server system 4 increases, the traffic (in the communication between servers) for making the information agree with each other increases, thereby deteriorating the throughput on the network 3.

[0056] Therefore, in the second embodiment, as in the configuration shown in Fig. 4, the private communication path 5 exclusive for information exchange between servers is provided between the respective file servers 10-1 to 10-n. The private communication path 5 is used for inter-server communication performed by the communication modules 413-1 to 413-n in the virtual distributed file modules 410-1 to 410-n in the virtual distributed file system 41, that is, for inter-server communication for making the contents of the mapping tables 14-1 to 14-n and the server information holding sections 15-1 to 15-n agree with each other.

[0057] In this embodiment, since the private communication path 5 is used, instead of the network 3, for the inter-server communication for making the contents of the mapping tables 14-1 to 14-n and the server information holding sections 15-1 to 15-n agree with each other, the load on the network 3 can be reduced.

[0058]

[Third Embodiment]

In the first and the second embodiments, a configuration example of the virtual distributed file server system that distributes a plurality of file servers so as to cooperate with each other is shown. The configuration shown in Figs. 1 and 4 referred to in the first and the second embodiments is a static example with a specific number of servers. However, it is desired that the number of servers can be changed.

[0059] Therefore, the third embodiment of the present invention in which the number of servers constituting the virtual distributed file server system can be changed will be explained, with reference to the accompanying drawings. Fig. 5 is a block diagram illustrating the configuration of a computer network system to which the virtual distributed file server system according to the third embodiment of the present invention is applied, wherein like reference signs refer to like parts in Fig. 4.

[0060] First, it is assumed that a new file server 10-(n+1) is added, as shown in Fig. 5(a), to the virtual distributed file server system 4 shown in Fig. 4, that is, the virtual distributed file server system 4 including n file servers 10-1 to 10-n.

[0061] In this case, a virtual distributed file module 410-(n+1) on the virtual distributed file system 41 distributed also on the added file server 10-(n+1) locks the update of the entry information in the mapping tables 14-1 to 14-n, and the server information (including the resource information and the load status of each server) in the server information holding sections 15-1 to 15-n, as shown by reference sign A1 in Fig. 5(a), with respect to the file servers 10-1 to 10-n already constituting the virtual distributed file server system 4, by inter-server communication (for example, via a private communication path (not shown)).

[0062] The module 410-(n+1) on the added server 10-(n+1) copies the whole information in the mapping table 14-1 and the server information holding section 15-1,

as shown by reference sign A2 in Fig. 5(b), to the mapping table 14-(n+1) and the server information holding section 15-(n+1) in the own server, from any one server of other file servers 10-1 to 10-n, for example, from the file server 10-1, by inter-server communication.

5 [0063] The module 410-(n+1) on the added server 10-(n+1) adds the server information indicating the resource and load status of the own server to the copied server information holding sections 15-(n+1), as shown by reference sign A3 in Fig. 5(c).

[0064] Thereafter, the module 410-(n+1) on the added server 10-(n+1) issues a
10 request for making the server information agree with each other with respect to all other file servers 10-1 to 10-n, as shown by reference sign A4 in Fig. 5(d), by inter-server communication, and then releases the lock.

[0065] By the series of operation, a new server (file server 10-(n+1)) can be added dynamically with respect to the virtual distributed file server system 4 already
15 constructed. In this case, for example, if information such as how to distribute the new resource with respect to the volume configuration of the current virtual distributed file server system 4 is added, the volume can be selectively expanded, according to need.

[0066]

[Fourth Embodiment]

20 Fig. 6 is a block diagram illustrating the configuration of a computer network system to which the virtual distributed file server system according to the fourth embodiment of the present invention is applied, wherein like reference signs refer to like parts in Fig. 4.

[0067] In Fig. 6, reference sign 6 denotes a virtual distributed file server system
25 corresponding to the virtual distributed file server system 4 in Fig. 4. The

characteristic point of the virtual distributed file server system 6 is that it includes each file load status information holding sections 16-1 to 16-n that hold the information of load status for each file (load status information for each file) held (in the storage devices 12-1 to 12-n) of the own server, in the file servers 10-1 to 10-n constituting the system 6. Therefore, the function of the virtual distributed file system 61, being the center of the virtual distributed file server system 6, is partially different from that of the virtual distributed file system 41 in Fig. 4. However, for the convenience sake, like reference signs (410-1 to 410-n) are used for the virtual distributed file modules for each file server 10-1 to 10-n in the virtual distributed file system 61. In Fig. 6, components (virtual distributed file interfaces, local file interfaces, and communication modules) in the virtual distributed file modules 410-1 to 410-n, and the mapping tables and the server information holding sections on the file servers 10-1 to 10-n are omitted.

[0068] The each file load status information holding section 16-i ($i = 1$ to n) has the data structure shown in Fig. 7(a), and holds the load status information for each file including the information indicating the load status for each file in the own server, the information indicating the attribute of the file (file attribute), and a replication flag. The file attribute indicates whether the corresponding file is the original or a replica (copy). When the corresponding file is the original, the replication flag indicates whether a replica thereof has already been generated on the other server side, that is, whether the replication has been already made.

[0069] In Fig. 6 is shown the situation in which a replica 612 of an optional file 611 in the storage device 12-1 included in the file server 10-1 is held by replication B1 in the storage device 12-n included in the file server 10-n.

[0070] The operation of the configuration shown in Fig. 6 will be explained. The virtual distributed file modules 410-1 to 410-n on the virtual distributed file system 61

refer to the each file load status information holding sections 16-1 to 16-n, for example, regularly. When having detected that a file having a load exceeding a first threshold exists in the files held by (the storage devices 12-1 to 12-n in) the own server, from the load status information for each file held by the holding sections 16-1 to 16-n, the
5 modules 410-1 to 410-n perform the replication operation for generating a replica of the corresponding file asynchronously with respect to one of the other servers, for example, via the private communication path 5 by inter-server communication. Here, the load status of the file is the sum total of the number of requests in a request queue with respect to the file, or the size indicated by the requests in the waiting state of the file,
10 and is updated every time a request is accepted and every time the request is processed. For the server to be replicated, for example, a server having the lowest load may be selected based on the server information held in the server information holding section (not shown).

[0071] When having performed the replication operation, the virtual distributed file
15 modules 410-1 to 410-n sets the replication flag in the load status information of the corresponding file held by the each file load status information holding sections 16-1 to 16-n in the own server to a replication notification state. The virtual distributed file module of the server to be replicated adds the load status information of the corresponding replica to the each file load status information holding section of the own
20 server.

[0072] Here, as shown in Fig. 6, it is assumed that replication B1 of the file 611 held by the file server 10-1 is performed with respect to the file server 10-n via the private communication path 5, and a replica 612 thereof is held in the storage device 12-n of the file server 10-n. In this case, the replication flag in the load status information of
25 the file 612 held by the each file load status information holding section 16-1 in the file

server 10-1 is set to the state indicating that replication has been finished. New load status information relating to the replica 612 of the file 611 is added to the each file load status information holding section 16-n in the file server 10-n. The file attribute in the load status information indicates that the corresponding file is a replica (612) (of the file 611).

[0073] Thereafter, when there is a new read request of the file 611 from the client 2, (the virtual distributed file modules 410-1 in) the file server 10-1 holding the file 611 checks if the load of the file 611 exceeds a second threshold (second threshold < first threshold), and if the load exceeds the second threshold, the file server 10-1 does not respond to the request from the client 2. In this case, the file server 10-n having accepted the replication responds to the request from the client 2. Here, the file server 10-n does not have to consider whether the file server 10-1 responds thereto, and may respond to the client 2, so long as it has the replica 612 of the requested file 611.

[0074] In this manner, since the file server 10-n processes the new read request with respect to the file 611 from the client 2 by using the replica 612, in the file server 10-1 holding the file 611, the processing with respect to the read request of the file 611, which has been accepted before, progresses, and the load of the file 611 becomes not larger than the second threshold. Then, the virtual distributed file module 410-1 on the file server 10-1 transmits a request for deleting the replica 612 of the file 611 to the virtual distributed file module 410-n on the file server 10-n, for example, by inter-server communication via the private communication path 5.

[0075] The virtual distributed file module 410-n on the file server 10-n having received the request performs processing only for the already accepted request by using the replica 612, and thereafter, deletes the replica 612 and the corresponding load status

information. On the other hand, when there is a new read request with respect to the file 611 from the client 2, the virtual distributed file module 410-1 on the file server 10-1 responds thereto.

[0076] Since the read request with respect to the file 611 is accepted by the file server
5 10-n, according to the replication of the file 611 from the file server 10-1 to the file server 10-n, the load of the replica 612 of the file 611 in the file server 10-n may exceed the first threshold, before the load of the file 611 in the file server 10-1 becomes not larger than the second threshold.

[0077] Therefore, in such a case, the file server 10-n generates the next generation
10 replica in another server by using the replica 612, that is, performs replication of the replica, and allows the server to handle the read request with respect to the file 611. For this purpose, as shown in Fig. 7(b), the load status information for each file held by the each file load status information holding section 16-i ($i = 1$ to n) may include the generation information of the file, in addition to the load status and the replication flag
15 as shown in Fig. 7(a).

[0078] In this case, at a point in time when the load of the replica of a certain generation becomes not larger than the second threshold, the server having the replica can perform management such that the next generation replica held by a server, to which the replication has been performed, is deleted. At this time, when the server, to
20 which deletion of the replica is requested, has generated the next generation replica in another server, the next generation replica can be also deleted. As for at least the load status information of the file relating to the replica, with regard to the same file (including the replica), a server holding the file having the lowest load may respond to the read request with respect to the file, by making the contents of information agree
25 with each other between respective servers, as with the server information described

above.

[0079] Recently, there is an increase in data, which is basically for readout, and has a relatively large size and requires a certain degree of response (band guarantee according to a case), such as streaming data for video and audio, or contents of the World Wide Web (WWW). For these data, since accesses may be concentrated on a specific data (file), securing the response may be difficult. The configuration in Fig. 6 assumes such a situation, and when accesses are concentrated on a specific file, replication of the file is automatically performed, so that the access to the file can be distributed. This configuration can be used not only for load balancing, but also for backup of a file of great importance.

[0080]

[Fifth Embodiment]

Fig. 8 is a block diagram illustrating the configuration of a computer network system to which the virtual distributed file server system according to the fifth embodiment of the present invention is applied, wherein like reference signs refer to like parts in Fig. 4.

[0081] In Fig. 8, reference sign 8 denotes a virtual distributed file server system corresponding to the virtual distributed file server system 4 in Fig. 4. The characteristic point of this virtual distributed file server system 8 is that a network configuration is applied such that the file servers 10-1 to 10-n and the storage devices 12-1 to 12-n are connected with each other, for example, by a fiber channel arbitrated loop (FC-AL) 80, so that the respective file servers 10-1 to 10-n (as a host) can share the storage devices 12-1 to 12-n (as a target) (that is, multi-hosting is possible). Here, it should be noted that the private communication path 5 is not provided, different from the configuration shown in Fig. 4.

[0082] In the configuration shown in Fig. 8, the inter-server communication, performed via the private communication path 5 (by the communication modules 413-1 to 413-n in the virtual distributed file modules 410-1 to 410-n) in the configuration of Fig. 4, may be performed via the network 3 as in the configuration shown in Fig. 1 (Fig. 8 illustrates this state). The inter-server communication may be performed on the FC-AL 80 via an interface for connecting the storages in the file servers 10-1 to 10-n. In this case, the load on the network 3 can be reduced, in the same manner when the private communication path 5 is used.

[0083] According to the configuration of Fig. 8, since the storage devices 12-1 to 12-n can be directly seen from all file servers 10-1 to 10-n, the replication operation and load balancing as described in the fourth embodiment can be easily performed, by providing the each file load status information holding sections 16-1 to 16-n shown in Fig. 6 to the respective servers 10-1 to 10-n. The network (interface) capable of multi-hosting is not limited to the FC-AL 80, and may be a small computer system interface (SCSI) bus.

[0084]

[Effects of the Invention]

According to the present invention described above, clients can handle a plurality of file servers distributed on the network as a single server, and clients do not have to consider the number of servers and the connection state of the storage devices.

[Brief Description of the Drawings]

[Fig. 1] A block diagram illustrating the configuration of a computer network system to which a virtual distributed file server system according to the first embodiment of the

present invention is applied.

[Fig. 2] A diagram illustrating an example of data structure of a mapping table shown in Fig. 1.

[Fig. 3] A diagram illustrating an example of data structure of a server information holding section shown in Fig. 1.

[Fig. 4] A block diagram illustrating the configuration of a computer network system to which the virtual distributed file server system according to the second embodiment of the present invention is applied.

[Fig. 5] A block diagram illustrating the configuration of a computer network system to which the virtual distributed file server system according to the third embodiment of the present invention is applied.

[Fig. 6] A block diagram illustrating the configuration of a computer network system to which the virtual distributed file server system according to the fourth embodiment of the present invention is applied.

[Fig. 7] A diagram illustrating an example of data structure of an each file load status information holding section shown in Fig. 6.

[Fig. 8] A timing chart for explaining the operation of the embodiment.

[Description of Signs]

1, 4, 6, 8 ... Virtual distributed file server system

20 2 ... Client

3 ... Network

5 ... Private communication path

10-1 to 10-n ... File server

11, 41, 61 ... Virtual distributed file system

25 12-1 to 12-n ... Storage device

- 13-1 to 13-n ... Local file system
- 14-1 to 14-n ... Mapping table
- 15-1 to 15-n ... Server information holding section
- 16-1 to 16-n ... Each file load status information holding section
- 5 80 ... FC-AL (interface capable of multi-hosting)
 - 110-1 to 110-n, 410-1 to 410-n ... Virtual distributed file module (Management module)
 - 111-1 to 111-n ... Virtual distributed file interface
 - 112-1 to 112-n ... Local file interface
 - 113-1 to 113-n ... Communication module
- 10 611 ... File
 - 612 ... Replica (of the file 611)

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号
特開2001-51890
(P2001-51890A)

(43) 公開日 平成13年2月23日 (2001.2.23)

(51) Int.Cl. ⁷	識別記号	F I	テーマコード [*] (参考)
G 0 6 F 12/00	5 4 5	G 0 6 F 12/00	5 4 5 A 5 B 0 8 2
13/00	3 5 1	13/00	3 5 1 E 5 B 0 8 9

審査請求 未請求 請求項の数 7 O L (全 14 頁)

(21) 出願番号 特願平11-226494
(22) 出願日 平成11年8月10日 (1999.8.10)

(71) 出願人 000003078
株式会社東芝
神奈川県川崎市幸区堀川町72番地
(72) 発明者 内堀 郁夫
東京都府中市東芝町1番地 株式会社東芝
府中工場内
(72) 発明者 高桑 正幸
東京都府中市東芝町1番地 株式会社東芝
府中工場内
(74) 代理人 100058479
弁理士 鈴江 武彦 (外6名)

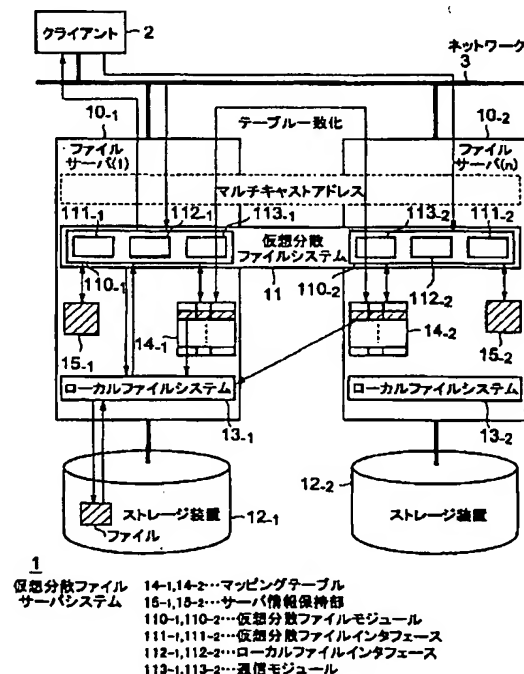
最終頁に続く

(54) 【発明の名称】 仮想分散ファイルサーバシステム

(57) 【要約】

【課題】 ネットワーク上に分散した複数のファイルサーバの台数やストレージ装置の接続状態をクライアントに意識させないで済むようにする。

【解決手段】 ネットワーク3上に分散したサーバ10-1, 10-2を備え、各サーバには、仮想分散ファイルシステム11が分散して実装されている。このシステム11を構成する、サーバ10-1, 10-2上のモジュール110-1, 110-2は、クライアント2からマルチキャストされたファイル操作要求を受け取ると、仮想分散ファイルシステム11と全ローカルファイルシステム13-1, 13-2とのマッピングテーブル14-1, 14-2または全サーバのサーバ情報を保持するサーバ情報保持部15-1, 15-2をもとに、自サーバが上記要求を処理可能な最適なサーバであるか否かを判断し、その判断結果に基づいて要求されたファイル操作を対応するサーバのローカルファイルシステムにより行わせる。



【特許請求の範囲】

【請求項 1】 マルチキャスト可能なネットワーク上に分散した複数のファイルサーバを備えた仮想分散ファイルサーバシステムであって、

前記各ファイルサーバに分散して実装され、全ファイルサーバのファイルを統合的に管理する、実際のストレージ構成には非依存の仮想分散ファイルシステムと、
前記各ファイルサーバにそれぞれ独立して実装され、各サーバに固有のストレージ構成を管理するローカルファイルシステムと、

前記各ファイルサーバにそれぞれ設けられ、前記仮想分散ファイルシステムで統合的に管理される各ファイルについて、当該仮想分散ファイルサーバシステムとそのファイルを実際に管理する前記ローカルファイルシステムとの間のマッピングの情報を保持するマッピングテーブルとを具備し、

前記仮想分散ファイルシステムは、前記各ファイルサーバにそれぞれ設けられた管理モジュールから構成され、前記各管理モジュールは、クライアントからマルチキャストされたファイル操作要求を共通に受け取り、当該要求に応じて自サーバの前記マッピングテーブルを参照することで、自サーバが当該要求を処理可能な最適なサーバであるか否かを判断し、最適なサーバであると判断した場合だけ、要求されたファイル操作を対応するサーバの前記ローカルファイルシステムにより行わせるように構成されていることを特徴とする仮想分散ファイルサーバシステム。

【請求項 2】 マルチキャスト可能なネットワーク上に分散した複数のファイルサーバを備えた仮想分散ファイルサーバシステムであって、

前記各ファイルサーバに分散して実装され、全ファイルサーバのファイルを統合的に管理する、実際のストレージ構成には非依存の仮想分散ファイルシステムと、
前記各ファイルサーバにそれぞれ独立して実装され、各サーバに固有のストレージ構成を管理するローカルファイルシステムと、

前記各ファイルサーバにそれぞれ設けられ、前記仮想分散ファイルシステムで統合的に管理される各ファイルについて、当該仮想分散ファイルサーバシステムとそのファイルを実際に管理する前記ローカルファイルシステムとの間のマッピングの情報を保持するマッピングテーブルと、

前記各ファイルサーバにそれぞれ設けられ、全ての前記ファイルサーバについて、そのサーバのストレージ装置の空き容量を示す情報、及びそのサーバの負荷状況を示す情報の少なくとも一方を含むサーバ情報を保持するサーバ情報保持手段とを具備し、

前記仮想分散ファイルシステムは、前記各ファイルサーバにそれぞれ設けられた管理モジュールから構成され、前記各管理モジュールは、クライアントからマルチキャストされた当該ファイルの読み出し要求があった場合、当

ストされたファイル操作要求を共通に受け取り、当該要求に応じて自サーバの前記マッピングテーブルまたは前記サーバ情報保持手段を参照することで、自サーバが当該要求を処理可能な最適なサーバであるか否かを判断し、最適なサーバであると判断した場合だけ、要求されたファイル操作を対応するサーバの前記ローカルファイルシステムにより行わせるように構成されていることを特徴とする仮想分散ファイルサーバシステム。

【請求項 3】 前記管理モジュールは、前記ファイル操作要求がファイル読み出し要求またはファイル書き込み要求の場合には、自サーバの前記マッピングテーブルを参照し、該当するファイルが自サーバの前記ローカルファイルシステムの管理下にあるか否かにより、自サーバが前記要求を処理可能な最適なサーバであるか否かを判断することを特徴とする請求項 1 または請求項 2 記載の仮想分散ファイルサーバシステム。

【請求項 4】 前記管理モジュールは、前記ファイル操作要求がファイルの新規作成要求の場合には、自サーバの前記サーバ情報保持手段を参照し、全ての前記サーバの各々について、そのサーバのストレージ装置の空き容量、またはそのサーバの負荷状況を比較することで、自サーバが前記要求を処理可能な最適なサーバであるか否かを判断することを特徴とする請求項 2 記載の仮想分散ファイルサーバシステム。

【請求項 5】 前記管理モジュールは、全ての前記ファイルサーバの前記マッピングテーブルの内容を一致化するために、自サーバの前記マッピングテーブルの情報と他のサーバの前記マッピングテーブルの情報とをサーバ間通信により交換することを特徴とする請求項 1 記載の仮想分散ファイルサーバシステム。

【請求項 6】 前記管理モジュールは、全ての前記ファイルサーバの前記マッピングテーブルの内容を一致化するために、自サーバの前記マッピングテーブルの情報と他のサーバの前記マッピングテーブルの情報とをサーバ間通信により交換する一方、全ての前記ファイルサーバの前記サーバ情報保持手段の内容を一致化するために、自サーバの前記サーバ情報保持手段の情報と他のサーバの前記サーバ情報保持手段の情報とをサーバ間通信により交換することを特徴とする請求項 2 記載の仮想分散ファイルサーバシステム。

【請求項 7】 前記各ファイルサーバにそれぞれ設けられ、そのサーバの管理下にある各ファイル別の負荷状況を示す情報を保持するファイル別負荷状況情報保持手段を更に具備し、

前記管理モジュールは、自サーバの前記ファイル別負荷状況情報保持手段に保持されている情報から第 1 の閾値を超えた負荷のファイルを検出して、他の任意のファイルサーバに対してサーバ間通信により当該ファイルのレプリケーションを行い、クライアントからマルチキャストされた当該ファイルの読み出し要求があった場合、当

該要求に対する処理をレプリケーション側に任せるようにしたことを特徴とする請求項 1 または請求項 2 記載の仮想分散ファイルサーバシステム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、コンピュータ・ネットワークシステムにおけるファイルサーバシステムに係り、特にネットワーク上に接続された複数のファイルサーバを協調動作させて、クライアントからは単一のサーバとして機能させる仮想分散ファイルサーバシステムに関する。

【0002】

【従来の技術】今日のコンピュータ・ネットワークシステムにおいては、ネットワークに接続された異なるコンピュータ間でファイルを共有することが一般的に行われている。こうした環境下では、特定のコンピュータに大規模なストレージを接続して、ファイルサーバとして運用したり、最近ではNAS (Network Attached Storage) と呼ばれる、ファイルサーバ専用機を接続する等のシステム構成をとることが多い。

【0003】ファイルサーバを使用する環境（ファイルサーバシステム）では、サーバのストレージ容量が不足した場合には、サーバ側に物理的・性能的に拡張性があれば、新たにディスク装置等（のストレージ装置）を増設することで対処できる。このときクライアントからは、新たなボリュームをマウントして使用するという形態になる。また、サーバの拡張性が限界に達していれば、サーバ自体を増設することになる。このときクライアントからは、増設したサーバを意識した上で新たなボリュームをマウントして使用するという形態になる。

【0004】

【発明が解決しようとする課題】上記したコンピュータ・ネットワークシステムにおいてファイルサーバを利用してファイル共有を行う場合、クライアント側からは、ファイルサーバ側のボリューム構成がそのまま見えてしまうのが一般的である。例えばサーバ側でディスク装置を増設した場合には、クライアント側は新たなボリュームを認識した上で、マウントしなければならない。或いはサーバ自体を増設した場合には、増設したサーバの運用ポリシーを決定、もしくはシステム設定・管理等の煩雑な作業が発生する上、クライアント側でも、新たなサーバを認識した上で、新たなボリュームをマウントしなければならない。

【0005】このように従来のファイルサーバを用いたファイル共有システム（ファイルサーバシステム）では、ディスク装置（ストレージ装置）の増設、或いはサーバの増設が必要な場合、サーバ側、クライアント側のいずれにも、新たな設定・管理のために多大なコストが発生するという問題があった。更に、ストレージの利用形態によっては、特定のファイルシステムをそのまま容量

だけ拡張したい場合もあり、単にストレージ装置やサーバを増設するだけでは解決しないケースもあった。

【0006】本発明は上記事情を考慮してなされたものでその目的は、ネットワーク上に分散した複数のファイルサーバを、クライアントからは単一のサーバとして扱うことができ、サーバ台数やストレージ装置の接続状態をクライアントに意識させない仮想分散ファイルサーバシステムを提供することにある。

【0007】

【課題を解決するための手段】本発明は、マルチキャスト可能なネットワークに接続された複数のファイルサーバに分散して実装され、全ファイルサーバのファイルを統合的に管理する、実際のストレージ構成には非依存の仮想分散ファイルシステムと、各ファイルサーバにそれぞれ独立して実装され、各サーバに固有のストレージ構成を管理するローカルファイルシステムと、前記各ファイルサーバにそれぞれ設けられ、上記各ファイルについて、仮想分散ファイルサーバシステムとそのファイルを実際に管理するローカルファイルシステムとの間のマッピングの情報（例えば、仮想分散ファイルサーバシステムで管理され、クライアントから見える仮想的なパスと、ローカルファイルシステムで管理され、クライアントから見えない物理的な所在とを対応付けた情報）を保持するマッピングテーブルとを備えると共に、上記仮想分散ファイルシステムを、各ファイルサーバにそれぞれ設けられた管理モジュールであって、クライアントからマルチキャストされたファイル操作要求を共通に受け取り、当該要求に応じて自サーバのマッピングテーブルを参照することで、自サーバが当該要求を処理可能な最適なサーバであるか否かを判断し、最適なサーバであると判断した場合だけ、要求されたファイル操作を対応するサーバのローカルファイルシステムにより行わせる管理モジュールにより構成することを特徴とする。

【0008】ここで、各ファイルサーバ上に、全ファイルサーバについて、そのサーバのストレージ装置の空き容量を示す情報、及びそのサーバの負荷状況を示す情報の少なくとも一方を含むサーバ情報を保持するサーバ情報保持手段を更に設け、上記各管理モジュールでは、クライアントからマルチキャストされたファイル操作要求を受け取った場合に、当該要求に応じて自サーバのマッピングテーブルまたはサーバ情報保持手段を参照することで、自サーバが当該要求を処理可能な最適なサーバであるか否かを判断する構成としてもよい。

【0009】このような構成においては、クライアントから特定のファイルサーバを意識することなくマルチキャストされたファイル操作要求は、仮想分散ファイルサーバシステムを構成する各ファイルサーバ上の管理モジュールで共通に受け取られ、その要求に応じて対応するサーバ（自サーバ）のマッピングテーブルまたはサーバ情報保持手段が参照される。そして、この参照の結果、

自サーバが上記要求を処理可能な最適なサーバであるか否かが判断され、最適なサーバであると判断した唯一のサーバ（上の管理モジュール）だけが、要求されたファイル操作を自サーバのローカルファイルシステムにより行わせる。

【0010】このように、要求元のクライアントからは、ネットワーク上に分散した複数のファイルサーバを単一のサーバとして扱うことができ、サーバ台数やストレージ装置の接続状態を意識する必要がない。

【0011】ここで、上記管理モジュールで、自サーバが最適なサーバであるか否かを判断するためのアルゴリズムとして、以下の第1乃至第4のアルゴリズム（判断手法）のいずれかを適用するとよい。

【0012】第1のアルゴリズムは、ファイル操作要求がファイル読み出し要求またはファイル書き込み要求の場合に適用されるもので、自サーバのマッピングテーブルの情報に基づいて、該当するファイルが自サーバのローカルファイルシステムの管理下にあるか否かにより判断する手法である。

【0013】第2のアルゴリズムは、ファイル操作要求がファイルの新規作成要求の場合に適用されるもので、自サーバのサーバ情報保持手段の情報に基づいて、全てのサーバの各々について、そのサーバのストレージ装置の空き容量（空き記憶容量）、またはそのサーバの負荷状況を比較することで判断する（例えば、自サーバの空き容量が最も大きい場合、或いは自サーバの負荷が最も低い場合に上記最適サーバと判断する）手法である。

【0014】第3のアルゴリズムも、ファイル操作要求がファイルの新規作成要求の場合に適用されるもので、自サーバのマッピングテーブルの情報に基づいて、全てのサーバの各々について対応するストレージ装置上に確保可能な連続領域を求め、その連続領域のサイズを比較することで判断する（例えば、自サーバのストレージ装置上に確保可能な連続領域のサイズが最も大きい場合に上記最適サーバと判断する）手法である。

【0015】第4のアルゴリズムも、ファイル操作要求がファイルの新規作成要求の場合に適用されるもので、全てのサーバの各々について、そのサーバのストレージ装置の空き容量、そのサーバの負荷、及び当該ストレージ装置上に確保可能な連続領域の少なくとも2つを求め、その求めた少なくとも2つの情報を複合条件として比較することで判断する手法である。

【0016】以上の第1乃至第4のアルゴリズムのいずれか1つを適用することで、クライアントから特定のファイルサーバを意識することなくマルチキャストされたファイル操作要求を各サーバが共通に受け取っても、その要求されたファイル操作を行うのに最適なサーバであるか否かを、その都度相互に通信を行うことなく、そのサーバ自身で自律的に判断することができる。

【0017】ここで、上記各管理モジュールに、全ての

ファイルサーバのマッピングテーブルの内容を一致化するために、自サーバのマッピングテーブルの情報と他サーバのマッピングテーブルの情報とをサーバ間通信により交換する機能（通信モジュール）を持たせるとよい。また、マッピングテーブルに加えてサーバ情報保持手段を各サーバ上に備えた構成では、各管理モジュール（内の通信モジュール）に、全てのファイルサーバのサーバ情報保持手段の内容を一致化するために、自サーバのサーバ情報保持手段の情報と他のサーバのサーバ情報保持手段の情報とをサーバ間通信により交換する機能を更に持たせるとよい。

【0018】また、マッピングテーブルの一致化のためには、自サーバのローカルファイルシステムで実際に管理されるファイル構成が変更された場合に、その変更された情報（マッピング情報）をサーバ間通信により他の全サーバに送信するのが効率的である。同様に、サーバ情報保持手段の内容の一致化のためには、自サーバのサーバ情報を定期的に更新し、その都度、その更新されたサーバ情報をサーバ間通信により他の全サーバに送信するのが効率的である。

【0019】また本発明は、上記仮想分散ファイルサーバシステムにサーバが動的に増設された場合に、そのサーバの管理モジュールで以下の第1乃至第4の処理を行うようにしたことをも特徴とする。まず、第1の処理では、サーバ間通信により他の全てのサーバに対してマッピングテーブル及びサーバ情報保持手段の更新を禁止するロック設定を行い、次の第2の処理では、サーバ間通信により他の任意のサーバからマッピングテーブルサーバ情報保持手段の内容を自サーバにコピーし、次の第3の処理では、自サーバのサーバ情報保持手段に自サーバのサーバ情報を追加し、次の第4の処理では、サーバ間通信により自サーバのサーバ情報を他の全てのサーバのサーバ情報保持手段に反映させて全サーバのサーバ情報保持手段の一致化を図り、しかる後に上記ロック設定を解除する。

【0020】このようなサーバ増設時の一連の動作により、動的にサーバ台数を拡張できる。しかもクライアントは、サーバ台数の拡張を意識することなく、増設されたサーバを利用することができる。

【0021】また本発明は、各ファイルサーバに、そのサーバの管理下にある各ファイル別の負荷状況を示す情報を保持するファイル別負荷状況情報保持手段を付加し、各サーバの管理モジュールにおいて、自サーバのファイル別負荷状況情報保持手段に保持されている情報から第1の閾値を超えた負荷のファイルを検出して、他の任意のファイルサーバに対してサーバ間通信により当該ファイルのレプリケーションを行い、クライアントからマルチキャストされた当該ファイルの読み出し要求があった場合、当該要求に対する処理をレプリケーション側に任せるようにしたことをも特徴とする。

【0022】このような構成においては、自律的な負荷分散が可能となる。ここで、上記検出したファイルがレプリケーションされたファイルである場合にも、他の任意のファイルサーバに対してサーバ間通信により当該ファイルのレプリケーションを行い、クライアントからマルチキャストされた当該ファイルの読み出し要求があった場合に、当該要求に対する処理を新たなレプリケーション側に任せることにより、自律的な負荷分散がより広範囲に効果的に行える。また、他サーバへのレプリケーションの対象となったファイルの負荷が前記第1の閾値より低い第2の閾値以下となった場合に、そのレプリケーションを行ったサーバ上の管理モジュールから当該他サーバに対してサーバ間通信により対応するファイルの消去を要求し、クライアントからマルチキャストされた当該ファイルの読み出し要求があった場合、当該要求に対する処理を自身が行うことにより、動的な負荷分散が可能となる。

【0023】さて、上記一致化のためのサーバ間通信（マッピング情報またはサーバ情報の通信）、更には上記レプリケーションのためのサーバ間通信には、上記ネットワークを用いることが可能である。しかし、各ファイルサーバを相互接続する専用の通信路（プライベート通信路）を上記ネットワークから独立に設け、当該通信路を用いてサーバ間通信を行う構成とするとよい。この場合、サーバ間通信のためにネットワークのスループットが悪化するのを防止できる。

【0024】また本発明は、各ファイルサーバ及び当該サーバのストレージ装置を相互接続するマルチホストが可能インタフェースを更に備えると共に、上記各管理モジュールによる上記サーバ間通信を当該インタフェースを介して行うようにしたことをも特徴とする。

【0025】このような構成においては、上記各ストレージ装置を上記インタフェースによって各サーバ間で共有し、上記一致化のためのサーバ間通信、更には動的なファイルのレプリケーションのためのサーバ間通信が上記インタフェースを通して行われるため、自律的な負荷分散が効果的に実現される。

【0026】なお、本発明は方法に係る発明としても成立する。

【0027】

【発明の実施の形態】以下、本発明の実施の形態につき図面を参照して説明する。

【0028】〔第1の実施形態〕図1は本発明の第1の実施形態に係る仮想分散ファイルサーバシステムを適用するコンピュータ・ネットワークシステムの構成を示すブロック図である。

【0029】同図において、1は仮想分散ファイルサーバシステム、2は同ファイルサーバシステム1にファイルサービスを要求するクライアント（クライアントコンピュータ）である。仮想分散ファイルサーバシステム1

は、ネットワーク3上に分散配置された複数、例えば2台のファイルサーバ（サーバコンピュータ）10-1、10-2を用いて実現される。なお、図ではクライアント2は便宜的に1台だけが示されているが、複数存在するのが一般的である。

【0030】11は仮想分散ファイルサーバシステム1の中心をなす仮想分散ファイルシステムであり、各ファイルサーバ10-1、10-2に分散して実装されている。この仮想分散ファイルシステム11は、全ファイルサーバ10-1、10-2のファイルを統合的に管理し、ファイルサーバ10-1、10-2それぞれの実際のボリューム構成（ストレージ構成）には依存しない、仮想的なファイルシステムをクライアント2に対して提供するものである。

【0031】仮想分散ファイルシステム11は、各ファイルサーバ10-1、10-2に分散して実装された仮想分散ファイルモジュール110-1、110-2を有している。仮想分散ファイルモジュール110-1、110-2は、ファイルサーバ10-1、10-2上でクライアント2からの要求を分散して処理しつつ、クライアント2に対しては仮想的に1つのファイルシステムとして見せるための管理モジュールである。仮想分散ファイルモジュール110-1、110-2は、当該モジュール110-1、110-2の中心をなし、クライアント2からの要求を処理する仮想分散ファイルインタフェース111-1、111-2と、後述するローカルファイルシステム13-1、13-2とのインタフェース（ローカルファイルインタフェース）112-1、112-2と、当該モジュール110-1、110-2間での通信（後述するマッピングテーブル14-1、14-2の情報、及びサーバ情報の一致化のための通信に代表されるサーバ間通信）を行う通信モジュール113-1、113-2を持つ。

【0032】ファイルサーバ10-1、10-2はネットワーク3を介してクライアント2と接続されている。ファイルサーバ10-1、10-2は、仮想分散ファイルシステム11の他に、それぞれ当該サーバ10-1、10-2に接続されたディスク装置などのストレージ装置12-1、12-2（実際のストレージ構成）を管理するローカルなファイルシステム（ローカルファイルシステム）13-1、13-2を実装している。

【0033】ファイルサーバ10-1、10-2には、仮想分散ファイルシステム11とローカルファイルシステム13-1、13-2とを対応付ける同一内容のマッピングテーブル14-1、14-2が設けられている。このテーブル14-i（i=1、2）のデータ構造を図2に示す。

【0034】テーブル14-iの各エントリは、仮想分散ファイルシステム1が管理するファイルのファイル名の登録フィールド（ファイル名フィールド）141と、当該ファイルの仮想分散ファイルシステム1上の論理的な所在を表す（クライアント2から見える）パス（仮想パ

ス)の登録フィールド(仮想パスフィールド)142、当該ファイルのストレージ上の物理的な所在位置を表す(クライアント2から見えない)所在位置情報の登録フィールド(所在位置情報フィールド)143、当該ファイルへのアクセス権(許可/禁止)を管理するためのパーミッション情報の登録フィールド(パーミッション情報フィールド)144、及びその他の各種属性の登録フィールド145を有している。

【0035】仮想分散ファイルシステム11(上の仮想分散ファイルモジュール110-i)は、このようなデータ構造のマッピングテーブル14-iを参照することにより、例えばあるファイルがファイルサーバ10-1、10-2のいずれにあるか等の所在情報を得ることができる。他、パーミッション等、必要に応じてファイルの属性を得ることができる。

【0036】ファイルサーバ10-1、10-2には更に、同一内容のサーバ情報保持部15-1、15-2が設けられている。サーバ情報保持部15-i(i=1,2)は、図3に示すように、仮想分散ファイルサーバシステム1を構成する全てのファイルサーバ10-1、10-2の(ストレージ装置12-1、12-2の)空き記憶容量を示す情報(リソース情報)、及び負荷状況を示す情報を含むサーバ情報を保持するのに用いられる。

【0037】次に図1の構成の動作を説明する。本実施形態では、クライアント2からは、各ファイルサーバ10-1、10-2のローカルファイルシステム13-1、13-2ではなくて、仮想分散ファイルシステム11がマウントされているように見えている。そこでクライアント2は、何らかのファイル操作要求が発生した場合、仮想分散ファイルシステム11が実装されている全ファイルサーバ10-1、10-2に対して同一の要求を発行する。この場合、例えばIP(Internet Protocol)マルチキャストを使用する等の手法によれば、クライアント2側はファイルサーバの台数を意識することなく要求の発行が可能である。

【0038】ファイルサーバ10-1、10-2は、クライアント2からの要求を受け取ると、当該要求を仮想分散ファイルシステム11内の自サーバに対応した仮想分散ファイルモジュール110-1、110-2に渡す。すると、モジュール110-1、110-2(内の仮想分散ファイルインタフェース111-1、111-2)は、その要求がファイルの読み出し要求もしくは書き込み(更新)要求であるか、または新規ファイルの作成要求もしくはディレクトリの作成要求であるか、その要求種別を判別する。

【0039】ここで、クライアント2からのファイル操作要求がファイルの読み出し要求もしくは書き込み要求であるものとする。この要求には、要求の対象となるファイルのファイル名と、当該ファイルの仮想分散ファイルシステム11上のパス(仮想パス)が付されている。

【0040】仮想分散ファイルモジュール110-1、110-2(内の仮想分散ファイルインタフェース111-1、111-2)は、クライアント2からのファイル操作要求がファイルの読み出し要求もしくは書き込み要求の場合、要求されたファイルのファイル名及び仮想パスにより自サーバ内のマッピングテーブル14-1、14-2を参照し、当該ファイル名及び仮想パスを持つテーブル14-1、14-2内エントリ中の所在位置情報フィールド143の登録情報から、操作要求のあったファイルが自サーバ内(自サーバに接続されたストレージ装置12-1、12-2)に保持されているか否かを調べる。

【0041】もし、要求されたファイルが自サーバ内に保持されている場合には、仮想分散ファイルモジュール110-1、110-2(内の仮想分散ファイルインタフェース111-1、111-2)は、ローカルファイルインタフェース112-1、112-2により自サーバ内のローカルファイルシステム13-1、13-2を介して実際のファイルにアクセスし、クライアント2に応答を返す。一方、操作要求のあったファイルが自サーバ内になかった場合には、他のサーバが応答するものと見なしに応答しない。

【0042】これに対し、クライアント2からの要求が新規ファイルの作成、或いはディレクトリの作成であった場合には、仮想分散ファイルモジュール110-1、110-2は、(マッピングテーブル14-1、14-2ではなくて)サーバ情報保持部15-1、15-2を参照する。そして、サーバ情報保持部15-1、15-2に保持されている全サーバのサーバ情報をもとに、所定のアルゴリズムに従い、いずれか1つのサーバ10-i(iは1または2)上の仮想分散ファイルモジュール110-iだけが、仮想分散ファイルインタフェース111-iによりクライアント2からの要求を受け付ける。具体的には、サーバ10-i上の仮想分散ファイルモジュール110-iは、全サーバのサーバ情報の示す空き記憶容量を比較し、自サーバ10-i(のストレージ装置12-i)の空き記憶容量が最も大きいと判定できる場合に、クライアント2からの要求を受け付けるものとする。この場合、必ずしもサーバ情報中に負荷状況の情報を持たせる必要はない。

【0043】なお、全サーバのサーバ情報の示す負荷を比較し、自サーバの負荷が最も低い場合にクライアント2からの要求を受け付けるようにしてもよい。この場合、必ずしもサーバ情報中に対応するサーバの空き記憶容量の情報を持たせる必要はない。

【0044】この他に、マッピングテーブル14-iの各ファイル毎のマッピング情報から(つまり各ストレージ装置12-1、12-2の領域の使用状況から)、各ストレージ装置12-1、12-2上に確保可能な連続領域を求め、必要なサイズ以上の連続領域が確保でき、且つそのサイズが最も大きいストレージ装置が自サーバのストレージ装置12-iの場合に、クライアント2からの要求を

受け付けるようにしてもよい。この場合、サーバ情報保持部15-1、15-2は必ずしも必要でない。

【0045】更に、空き記憶容量と負荷状況と確保できる連続領域のサイズの少なくとも2つを条件（複合条件）として評価値を求め、自サーバが要求を受け付ける最適なサーバであるか否かを判断ようにしてもよい。

【0046】さて、ファイルサーバ10-i上の仮想分散ファイルモジュール110-iでは、仮想分散ファイルインタフェース111-iによりクライアント2からの要求を受け付けると、要求された新規ファイルの作成、或いはディレクトリの作成を、ローカルファイルインタフェース112-iを介してローカルファイルシステム13-iにより行い、マッピングテーブル14-1、14-2に該当するエントリ情報を登録する。

【0047】新規ファイルの作成、或いはディレクトリの作成が完了した後は、ファイルサーバ10-i上の仮想分散ファイルモジュール110-iでは、自サーバ上のマッピングテーブル14-iに登録した新たなエントリ情報を通信モジュール113-iによりネットワーク3を介して他の全てのサーバ10-j（jは1または2、但しj≠i）上の仮想分散ファイルモジュール110-jに送る。仮想分散ファイルモジュール110-j（内の仮想分散ファイルインタフェース111-j）は、仮想分散ファイルモジュール110-iから送られたマッピングテーブル14-iのエントリ情報を通信モジュール113-jを介して受け取る。そしてモジュール110-j（内のインタフェース111-j）は、受け取った他サーバ10-iのマッピングテーブル14-iのエントリ情報を自サーバ内のマッピングテーブル14-jに登録する。このように、ファイルサーバ10-1、110-2上の仮想分散ファイルモジュール110-1、110-2が相互にマッピングテーブル14-1、14-2の新規登録されたエントリ情報（更には更新されたエントリ情報）を交換し合うことで、当該マッピングテーブル14-1、14-2の内容の一致化を図ることができる。

【0048】また、ファイルサーバ10-1、10-2上の仮想分散ファイルモジュール110-1、110-2は、自サーバ上のサーバ情報保持部15-1、15-2に保持されている各サーバのサーバ情報のうち、自サーバのサーバ情報（空き記憶容量、及び負荷状況）を定期的に更新すると共に、その更新したサーバ情報を（通信モジュール113-1、113-2により）ネットワーク3を介して他の全てのサーバ（上の仮想分散ファイルモジュール110-2、110-1）に定期的に送ることで、各ファイルサーバ10-1、10-2のサーバ情報保持部15-1、15-2の内容の一致化を図る。つまり仮想分散ファイルモジュール110-1、110-2は定期的にサーバ情報を交換し合うことで一致化を図る。

【0049】以上の動作によって、ファイルサーバ10-1、10-2を自律的に分散・協調動作させることがで

き、クライアント2には実際にはファイルサーバが2台（複数台）あることを意識させずに、仮想的なファイルサーバを提供することができる。

【0050】なお、図1のシステムの例ではサーバが2台である場合について説明したが、サーバが3台以上であっても同様の仕組みによって、仮想的なファイルサーバを提供することができる。

【0051】〔第2の実施形態〕図4は本発明の第2の実施形態に係る仮想分散ファイルサーバシステムを適用するコンピュータ・ネットワークシステムの構成を示すブロック図であり、図1と同一部分には同一符号を付してある。

【0052】図4において、仮想分散ファイルサーバシステム4の中心をなす仮想分散ファイルシステム41は、n台のファイルサーバ10-1～10-nに分散して実装されている。この仮想分散ファイルシステム41は、図1中の仮想分散ファイルシステム11と同様に、全ファイルサーバ10-1～10-nのファイルを統合的に管理し、各ファイルサーバ10-1～10-nそれぞれの実際のボリューム構成には依存しない、仮想的なファイルシステムをクライアント2に対して提供している。仮想分散ファイルシステム41は、各ファイルサーバ10-1～10-n上でクライアント2からの要求を処理する仮想分散ファイルモジュール410-1～410-nを有している。モジュール410-1～410-nは、図1中のモジュール110-1、110-2と同様の構成である、仮想分散ファイルインタフェース411-1～411-nと、ローカルファイルインタフェース412-1～412-nと、通信モジュール413-1～413-nとを持つ。但し、本実施形態における通信モジュール413-1～413-nは、図1中の通信モジュール113-1、113-2と異なり、後述するプライベート通信路5を介して通信を行うように構成されている。

【0053】ファイルサーバ10-1～10-nはネットワーク3を介してクライアント2と接続されている。ファイルサーバ10-1～10-nは、仮想分散ファイルシステム11の他に、それぞれ当該サーバ10-1～10-nに接続されたストレージ装置12-1～12-2を管理するローカルなファイルシステム（ローカルファイルシステム）13-1～13-2を実装している。ファイルサーバ10-1～10-nには、マッピングテーブル14-1～14-nと、サーバ情報保持部15-1～15-2とが設けられている。

【0054】図4の構成の仮想分散ファイルサーバシステム4の特徴は、図1の構成の仮想分散ファイルサーバシステム1と異なって、システムを構成するファイルサーバの台数がn台である点と、そのn台のファイルサーバ10-1～10-nがネットワーク3とは別のプライベート通信路5によっても相互接続されている点である。このプライベート通信路5は、例えばイーサネット、或いはファイバチャネル（Fibre Channel）等であるが、

物理層に関しては特定しない。またトポロジに関しても、図4の例ではバス型を想定しているが、ループやスイッチであってもよい。

【0055】図4の構成において、ファイルサーバ10-1～10-nが（仮想分散ファイルシステム41内の仮想分散ファイルモジュール410-1～410-nにより）分散・協調動作を行うためには、前記第1の実施形態での動作説明から類推されるように、マッピングテーブル14-1～14-n、及びサーバ情報保持部15-1～15-nの内容を、各サーバ10-1～10-n間で常に一致化しておく必要がある。しかし、サーバ10-1～10-n間の情報一致化を、前記第1の実施形態と同様にネットワーク3を介して行うのでは、仮想分散ファイルサーバシステム4を構成するファイルサーバの台数が増加した場合には、その情報一致化の（ためのサーバ間通信の）トラフィックが増加し、ネットワーク3上のスループットを悪化させることになる。

【0056】そこで本実施形態（第2の実施形態）では、図4の構成のように、各ファイルサーバ10-1～10-n間にサーバ間の情報交換専用のプライベート通信路5を設け、仮想分散ファイルシステム41内の仮想分散ファイルモジュール410-1～410-nで通信モジュール413-1～413-nにより行われるサーバ間通信に、即ちマッピングテーブル14-1～14-n、及びサーバ情報保持部15-1～15-nの内容を一致化するためのサーバ間通信に、この通信路5を使用するようにしている。

【0057】このように本実施形態では、マッピングテーブル14-1～14-n、及びサーバ情報保持部15-1～15-nの内容の一致化のためのサーバ間通信に、ネットワーク3でなくてプライベート通信路5を用いることにより、ネットワーク3の負荷の軽減を図ることができる。

【0058】〔第3の実施形態〕以上に述べた第1、第2の実施形態では、複数のファイルサーバを分散・協調動作させる仮想分散ファイルサーバシステムの構成例を示した。この第1、第2の実施形態で参照した図1、図4の構成は、特定のサーバ台数における静的な例である。しかし、サーバ台数については、変更可能な構成とすることが好ましい。

【0059】そこで、仮想分散ファイルサーバシステムを構成するサーバ台数を動的に拡張可能とした本発明の第3の実施形態について図面を参照して説明する。図5は本発明の第3の実施形態に係る仮想分散ファイルサーバシステムを適用するコンピュータ・ネットワークシステムの構成を示すブロック図であり、図4と同一部分には同一符号を付してある。

【0060】まず、図4に示した仮想分散ファイルサーバシステム4、即ちn台のファイルサーバ10-1～10-nで構成される仮想分散ファイルサーバシステム4に、図5（a）に示すように、新たなファイルサーバ10-

(n+1)を追加するものとする。

【0061】この場合、追加されたファイルサーバ10-(n+1)にも分散されている仮想分散ファイルシステム41上の仮想分散ファイルモジュール410-(n+1)は、既に仮想分散ファイルサーバシステム4を構成しているファイルサーバ10-1～10-nに対し、図5（a）において符号A1で示すように、（例えば図示せぬプライベート通信路を介しての）サーバ間通信により、マッピングテーブル14-1～14-nのエントリ情報及びサーバ情報保持部15-1～15-nのサーバ情報（各サーバのリソース情報及び負荷状況を含む）の更新をロックする。

【0062】その上で、追加されたサーバ10-(n+1)上のモジュール410-(n+1)は、他のファイルサーバ10-1～10-nのうちのいずれかのサーバ、例えばファイルサーバ10-1から、図5（b）において符号A2で示すように、マッピングテーブル14-1及びサーバ情報保持部15-1の全情報を、サーバ間通信により自サーバ内のマッピングテーブル14-(n+1)及びサーバ情報保持部15-(n+1)にコピーする。

【0063】次に、追加されたファイルサーバ10-(n+1)上のモジュール410-(n+1)は、コピー後のサーバ情報保持部15-(n+1)に対し、図5（c）において符号A3で示すように、自サーバのリソース及び負荷状況を示すサーバ情報を追加する。

【0064】しかる後にファイルサーバ10-(n+1)上のモジュール410-(n+1)は、図5（d）において符号A4で示すように、サーバ間通信により他の全ファイルサーバ10-1～10-nに対してサーバ情報の一致化要求を発行し、その後にロックを解除する。

【0065】以上の一連の動作により、既に構築されている仮想分散ファイルサーバシステム4に対して、動的に新たなサーバ（ファイルサーバ10-(n+1)）を追加することができる。この場合、例えば現在の仮想分散ファイルサーバシステム4のボリューム構成に対し、新規リソースをどのように振り分けるか、といった情報を付加すれば、必要に応じてボリュームを選択的に拡張することも可能である。

【0066】〔第4の実施形態〕図6は本発明の第4の実施形態に係る仮想分散ファイルサーバシステムを適用するコンピュータ・ネットワークシステムの構成を示すブロック図であり、図4と同一部分には同一符号を付してある。

【0067】図6において、6は図4中の仮想分散ファイルサーバシステム4に相当する仮想分散ファイルサーバシステムである。この仮想分散ファイルサーバシステム6の特徴は、当該システム6を構成するファイルサーバ10-1～10-n内に、自サーバ（のストレージ装置12-1～12-n）に保持されている各ファイルについての負荷状況の情報（ファイル別負荷状況情報）を保持するファイル別負荷状況情報保持部16-1～16-nを備えて

いる点にある。これに伴い、仮想分散ファイルサーバシステム6の中心をなす仮想分散ファイルシステム61の持つ機能も図4中の仮想分散ファイルシステム41とは一部異なる。但し、仮想分散ファイルシステム61内の各ファイルサーバ10-1~10-n毎の仮想分散ファイルモジュールには便宜的に図4と同一符号(410-1~410-n)を用いている。なお、図6では、モジュール410-1~410-n内の構成要素(仮想分散ファイルインタフェース、ローカルファイルインタフェース、通信モジュール)、及びファイルサーバ10-1~10-n上のマッピングテーブル、サーバ情報保持部は省略されている。

【0068】ファイル別負荷状況情報保持部16-i(i=1~n)は、図7(a)に示すデータ構造を持ち、自サーバ内のファイル毎の負荷状況を示す情報と、そのファイルの属性を示す情報(ファイル属性)と、レプリケーションフラグとを含むファイル別負荷状況情報を保持する。ファイル属性は、対応するファイルがオリジナルであるかレプリカ(複製)であるかを示す。また、レプリケーションフラグは、対応するフラグがオリジナルの場合、そのレプリカを他サーバ側に生成済みであるか否か、つまりレプリケーション済みであるか否かを示す。

【0069】図6には、ファイルサーバ10-nが持つストレージ装置12-n内に、ファイルサーバ10-1が持つストレージ装置12-1内の任意のファイル611のレプリカ612がレプリケーションB1により保持されている様子が示されている。

【0070】次に、図6の構成の動作を説明する。仮想分散ファイルシステム61上の各仮想分散ファイルモジュール410-1~410-nは、ファイル別負荷状況情報保持部16-1~16-nを例えば定期的に参照する。そしてモジュール410-1~410-nは、保持部16-1~16-nに保持されているファイル別の負荷状況情報から、自サーバ(のストレージ装置12-1~12-n)に保持されているファイルの中に、第1の閾値を超えた負荷のファイルが存在することを検出した場合、他サーバの1つに対して対応するファイルのレプリカを非同期に生成するレプリケーション動作を、例えばプライベート通信路5を介してのサーバ間通信により行う。ここでファイルの負荷状況は、当該ファイルへの要求の待ち行列(キュー)にある要求数、或いは当該ファイルの待ち状態にある要求の示すサイズの総和であり、要求を受け付ける毎と要求を処理し終える毎に更新される。また、レプリケーションの対象サーバには、図示せぬサーバ情報保持部に保持されているサーバ情報に基づいて、例えば負荷が最も低いサーバを選択すればよい。

【0071】仮想分散ファイルモジュール410-1~410-nはレプリケーション動作を行うと、自サーバのファイル別負荷状況情報保持部16-1~16-nに保持されている対応するファイルの負荷状況情報中のレプリケ

ションフラグをレプリケーション済みの通知状態にセットする。またレプリケーション動作の対象となったサーバの仮想分散ファイルモジュールは、自サーバのファイル別負荷状況情報保持部内に対応するレプリカの負荷状況情報を追加する。

【0072】ここでは、図6に示すように、ファイルサーバ10-1の保持するファイル611のレプリケーションB1がプライベート通信路5を介してファイルサーバ10-nに対して行われて、そのレプリカ612が当該ファイルサーバ10-nのストレージ装置12-nに保持されたものとする。この場合、ファイルサーバ10-1のファイル別負荷状況情報保持部16-1に保持されているファイル611の負荷状況情報中のレプリケーションフラグがレプリケーション済みを示す状態にセットされる。また、ファイルサーバ10-nのファイル別負荷状況情報保持部16-nには、ファイル611のレプリカ612についての新たな負荷状況情報が追加される。この負荷状況情報中のファイル属性は、対応するファイルが(ファイル611の)レプリカ(612)であることを示す。

【0073】以後、クライアント2からファイル611の新たな読み出し要求があった場合、当該ファイル611を保持するファイルサーバ10-1(の仮想分散ファイルモジュール410-1)は、当該ファイル611の負荷が第2の閾値(但し、第2の閾値<第1の閾値)を超えているか否かを調べ、超えているならば、クライアント2からの要求に応答しない。この場合、クライアント2からの要求に対しては、レプリケーションを受けたファイルサーバ10-nが応答する。ここでファイルサーバ10-nは、ファイルサーバ10-1が応答するか否かを考慮する必要はなく、要求されたファイル611のレプリカ612を有する限り、クライアント2に応答すればよい。

【0074】このように、クライアント2からのファイル611に対する新たな読み出し要求を、そのレプリカ612を用いてファイルサーバ10-nが処理することで、そのファイル611を保持するファイルサーバ10-1では、それ以前に受け付けた当該ファイル611に対する読み出し要求の処理が進み、当該ファイル611の負荷が上記第2の閾値以下となる。するとファイルサーバ10-1上の仮想分散ファイルモジュール410-1は、ファイルサーバ10-n上の仮想分散ファイルモジュール410-nに対して、ファイル611のレプリカ612を消去するための要求を例えばプライベート通信路5を介したサーバ間通信により送る。

【0075】この要求を受けたファイルサーバ10-n上の仮想分散ファイルモジュール410-nは、既に受け付け済みの要求に対してのみレプリカ612を用いて処理を行い、しかる後にレプリカ612と対応する負荷状況情報を消去する。一方、ファイルサーバ10-1上の仮想分散ファイルモジュール410-1は、クライアント2か

らのファイル 611 に対する新たな読み出し要求があれば、それに対して応答する。

【0076】ところで、ファイルサーバ 10-1 からファイルサーバ 10-n へのファイル 611 のレプリケーションにより、当該ファイル 611 に対する読み出し要求をファイルサーバ 10-n で受け付けるようになった結果、ファイルサーバ 10-1 におけるファイル 611 の負荷が第 2 の閾値以下となる前に、ファイルサーバ 10-n における当該ファイル 611 のレプリカ 612 の負荷が第 1 の閾値を超えることがあり得る。

【0077】そこで、このような場合、今度はファイルサーバ 10-n がレプリカ 612 を用いて次の世代のレプリカを他の 1 つのサーバに生成し、即ちレプリケーションのレプリケーションを行い、そのサーバでファイル 611 に対する読み出し要求を処理させればよい。そのためには、ファイル別負荷状況情報保持部 16-i (i=1~n) に保持されるファイル毎の負荷状況情報に、図 7 (a) に示したような負荷状況とレプリケーションフラグに加えて、図 7 (b) に示すように、ファイルの世代情報を持たせるとよい。

【0078】この場合、ある世代のレプリカの負荷が上記第 2 の閾値以下に下がった時点で、当該レプリカを持つサーバから、そのサーバによるレプリケーションの対象となったサーバの持つ次世代のレプリカを消去する等の制御を行うことができる。このとき、レプリカの消去が要求されたサーバが、別のサーバに対して更に次世代のレプリカを生成している場合、その更に次世代のレプリカを消去するとよい。この他に、少なくともレプリカに関連したファイルの負荷状況情報については、前記したサーバ情報と同様に、各サーバ間の一致化を図ることにより、同一ファイル (レプリカを含む) について、負荷が最も低いファイルを持つサーバが、当該ファイルに対する読み出し要求に応答するようにしてもよい。

【0079】最近では、ビデオ、オーディオ等のストリーミングデータや、或いは WWW (World Wide Web) のコンテンツ等、基本的には読み出しが主で、比較的サイズが大きく、ある程度のレスポンス (場合によっては帯域保証) が必要なデータが増加しつつある。しかもこうしたデータは、短期的に見て特定のデータ (ファイル) にアクセスが集中するケースが想定されるため、レスポンスを確保するのが困難な場合もある。以上に述べた図 6 の構成は、こうした状況を想定したもので、特定のファイルにアクセスが集中した場合に、自動的に当該ファイルのレプリケーションを行うことで、当該ファイルへのアクセスを分散させることができるようにしている。この構成は、単に負荷分散だけでなく、例えば重要性の高いファイルのバックアップに利用することも可能である。

【0080】〔第 5 の実施形態〕図 8 は本発明の第 5 の実施形態に係る仮想分散ファイルサーバシステムを適用

するコンピュータ・ネットワークシステムの構成を示すブロック図であり、図 4 と同一部分には同一符号を付してある。

【0081】図 8 において、8 は図 4 中の仮想分散ファイルサーバシステム 4 に相当する仮想分散ファイルサーバシステムである。この仮想分散ファイルサーバシステム 8 の特徴は、ファイルサーバ 10-1 ~ 10-n、及びストレージ装置 12-1 ~ 12-n が、例えば FC-A L (Fibre Channel Arbitrated Loop) 80 により相互接続され、(ホストとしての) 各ファイルサーバ 10-1 ~ 10-n から (ターゲットとしての) ストレージ装置 12-1 ~ 12-n の共有が可能 (つまりマルチホスト可能な) ネットワーク構成を適用している点にある。ここでは、図 4 の構成と異なって、プライベート通信路 5 を持たない点に注意されたい。

【0082】この図 8 の構成では、図 4 の構成において (仮想分散ファイルモジュール 410-1 ~ 410-n の通信モジュール 413-1 ~ 413-n により) プライベート通信路 5 を介して行われるサーバ間通信を、図 1 の構成と同様にネットワーク 3 を介して行えばよい (図は、この状態が示されている)。また、上記サーバ間通信を、ファイルサーバ 10-1 ~ 10-n のストレージ接続用のインタフェースを介して FC-A L 80 上で行うようにしてもよい。この場合、プライベート通信路 5 を用いたのと同様に、ネットワーク 3 の負荷を軽減できる。

【0083】図 8 の構成によれば、ストレージ装置 12-1 ~ 12-n が全てのファイルサーバ 10-1 ~ 10-n から直接に見えるので、各サーバ 10-1 ~ 10-n に図 6 中のファイル別負荷状況情報保持部 16-1 ~ 16-n を持たせることで、前記第 4 の実施形態で述べたようなレプリケーション動作や負荷分散を容易に行うことができる。なお、マルチホスト可能なネットワーク (インタフェース) は FC-A L 80 に限るものではなく、SCSI (Small Computer System Interface) バスであっても構わない。

【0084】

【発明の効果】以上詳述したように本発明によれば、ネットワーク上に分散した複数のファイルサーバを、クライアントからは単一のサーバとして扱うことができ、サーバ台数やストレージ装置の接続状態をクライアントに意識させることがない。

【0085】また本発明によれば、サーバを増設した場合、動的にボリュームを拡張することもできる。

【0086】更に本発明によれば、複数のサーバ間で自律的な負荷分散が実現できる。

【図面の簡単な説明】

【図 1】本発明の第 1 の実施形態に係る仮想分散ファイルサーバシステムを適用するコンピュータ・ネットワークシステムの構成を示すブロック図。

【図 2】図 1 中のマッピングテーブルのデータ構造例を

示す図。

【図3】図1中のサーバ情報保持部のデータ構造例を示す図。

【図4】本発明の第2の実施形態に係る仮想分散ファイルサーバシステムを適用するコンピュータ・ネットワークシステムの構成を示すブロック図。

【図5】本発明の第3の実施形態に係る仮想分散ファイルサーバシステムを適用するコンピュータ・ネットワークシステムの構成を示すブロック図。

【図6】本発明の第4の実施形態に係る仮想分散ファイルサーバシステムを適用するコンピュータ・ネットワークシステムの構成を示すブロック図。

【図7】図6中のファイル別負荷状況情報保持部のデータ構造例を示す図。

【図8】同実施形態の動作を説明するタイミングチャート。

【符号の説明】

- 1, 4, 6, 8…仮想分散ファイルサーバシステム
2…クライアント

3…ネットワーク

5…プライベート通信路

10-1～10-n…ファイルサーバ

11, 41, 61…仮想分散ファイルシステム

12-1～12-n…ストレージ装置

13-1～13-n…ローカルファイルシステム

14-1～14-n…マッピングテーブル

15-1～15-n…サーバ情報保持部

16-1～16-n…ファイル別負荷状況情報保持部

10 80…F C - A L (マルチホスト可能なインタフェース)

110-1～110-n, 410-1～410-n…仮想分散ファイルモジュール (管理モジュール)

111-1～111-n…仮想分散ファイルインタフェース

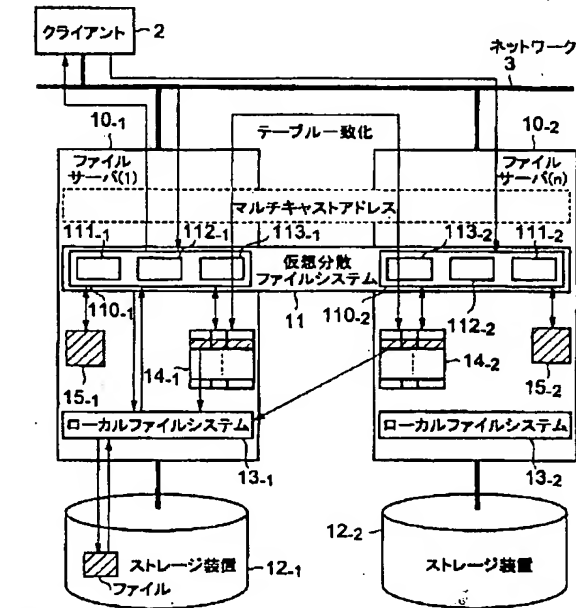
112-1～112-n…ローカルファイルインタフェース

113-1～113-n…通信モジュール

611…ファイル

612… (ファイル611の) レプリカ

【図1】



- 1
仮想分散ファイル
サーバシステム
- 14-1, 14-2…マッピングテーブル
15-1, 15-2…サーバ情報保持部
110-1, 110-2…仮想分散ファイルモジュール
111-1, 111-2…仮想分散ファイルインタフェース
112-1, 112-2…ローカルファイルインタフェース
113-1, 113-2…通信モジュール

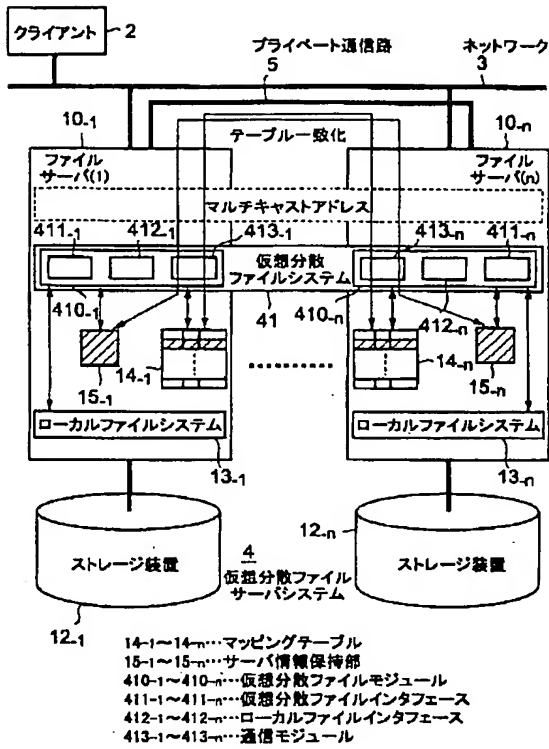
【図2】

141	142	143	144	145	145
ファイル名	仮想パス	所在位置情報	パーミッション情報	その他の属性	その他の属性
ファイル名	仮想パス	所在位置情報	パーミッション情報	その他の属性	その他の属性
ファイル名	仮想パス	所在位置情報	パーミッション情報	その他の属性	その他の属性

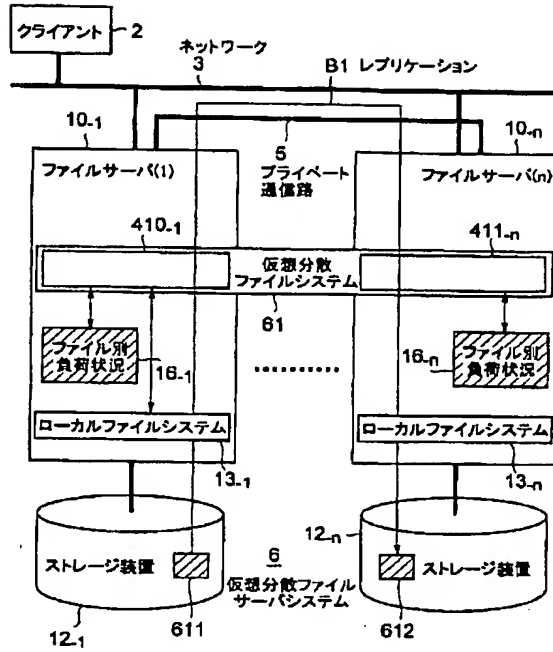
【図3】

サーバID	リソース情報	負荷状況
サーバ1		
サーバ2		

【図 4】

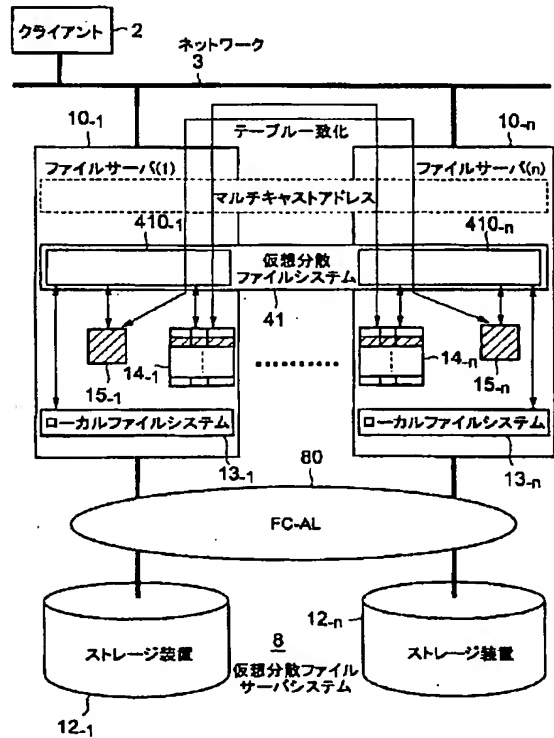
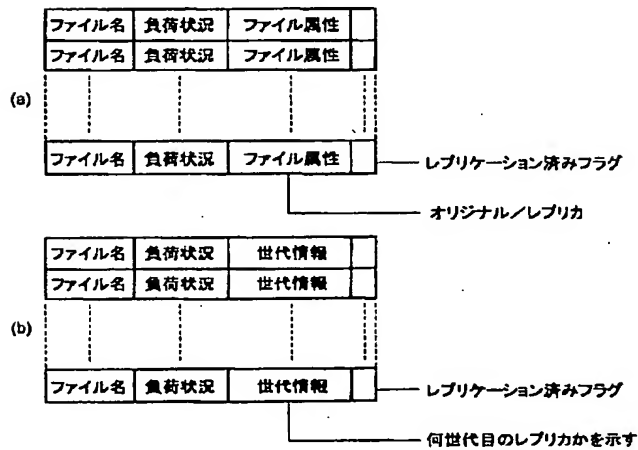


【図 6】

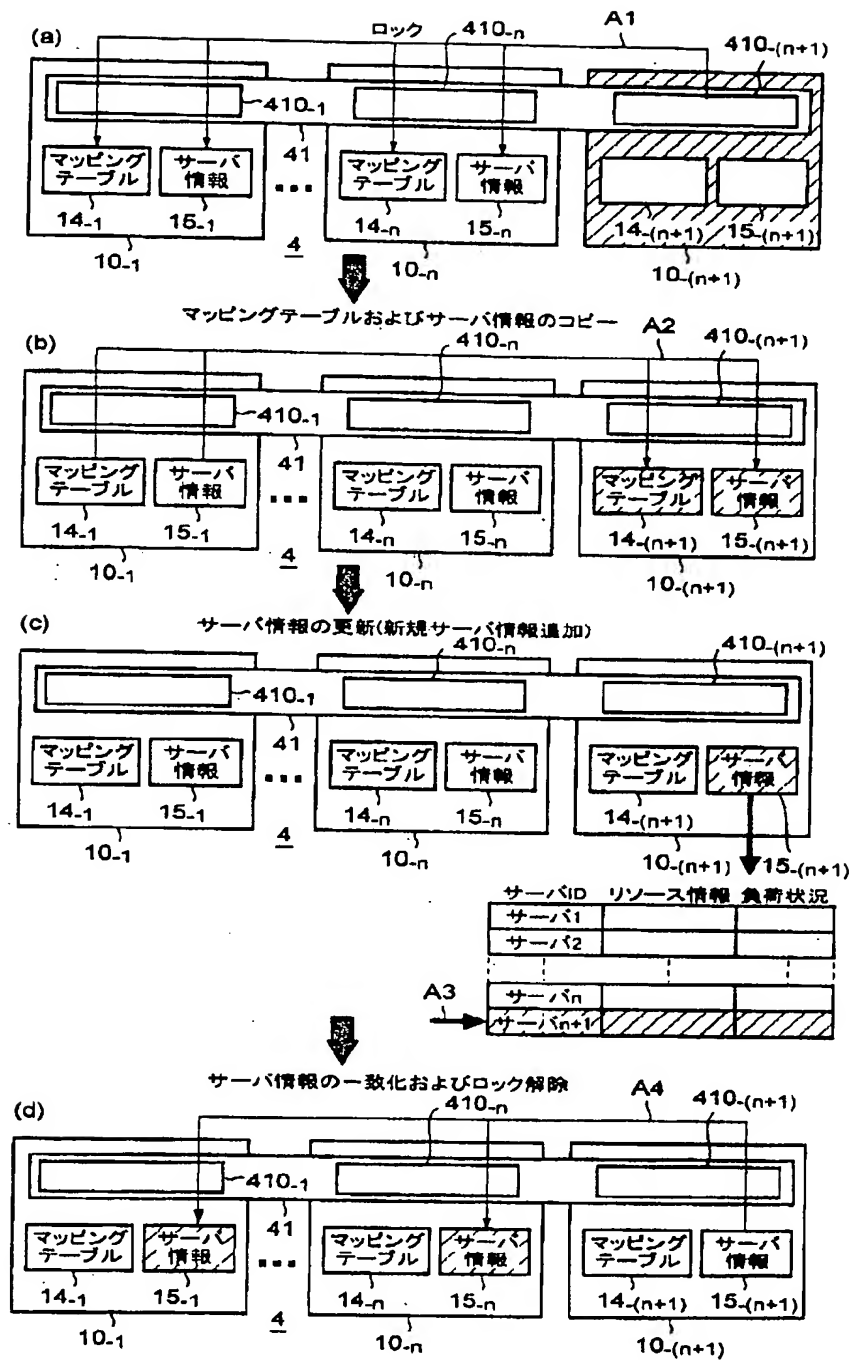


【図 8】

【図 7】



【図5】



フロントページの続き

Fターム(参考) 5B082 CA18 EA07 HA03 HA05 HA08
HA09
5B089 GA12 JA11 JB15 KA00 KC15
KC28 KE07